# Modeling Worldwide Gross Box Office Collection for Films 2012–2021

Rachel Goodridge

# Model Motivation

What factors can be used to predict how much a movie will make in theaters worldwide?

# Tools and Methods

1. Scrape target (Worldwide Gross) and features from films listed in 2012-2021 on Box Office Mojo
2. Investigate dependent and independent variables
3. Build models using the LinearRegression function from the sklearn package
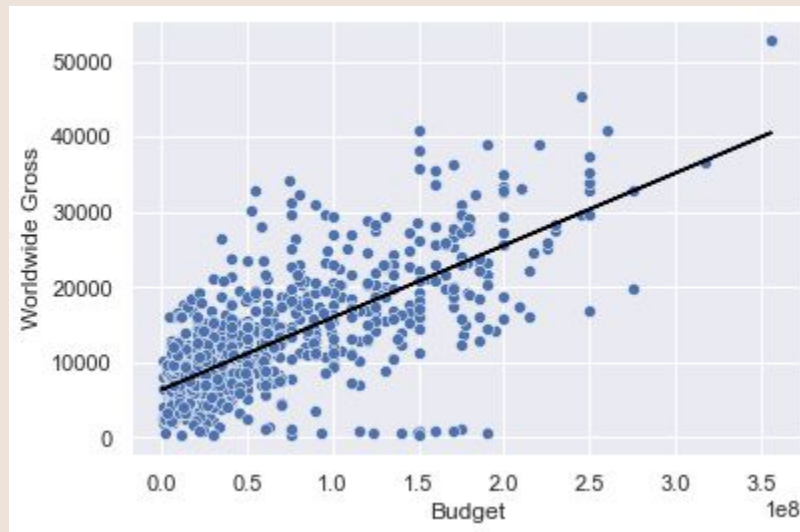4. Choose the best model based on scores from cross-validation

# Feature Investigation

- Budget
- Running Time
- Number of Theaters
- Distributor
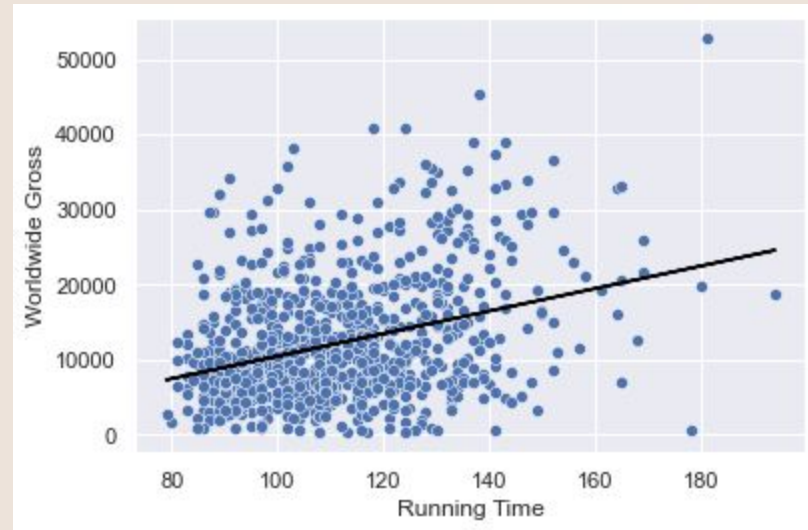- Genre
- Release Year

# Feature Investigation

- **Budget**
- Running Time
- Number of Theaters
- Distributor
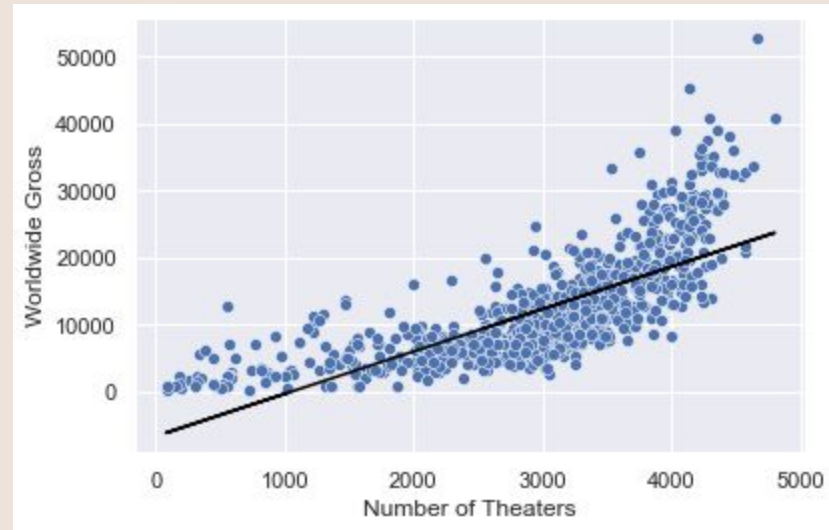- Genre
- Release Year

# Feature Investigation

- Budget
- **Running Time**
- Number of Theaters
- Distributor
- Genre
- Release Year

# Feature Investigation

- Budget
- Running Time
- **Number of Theaters**
- Distributor
- Genre
- Release Year

# Feature Investigation

- Budget
- Running Time
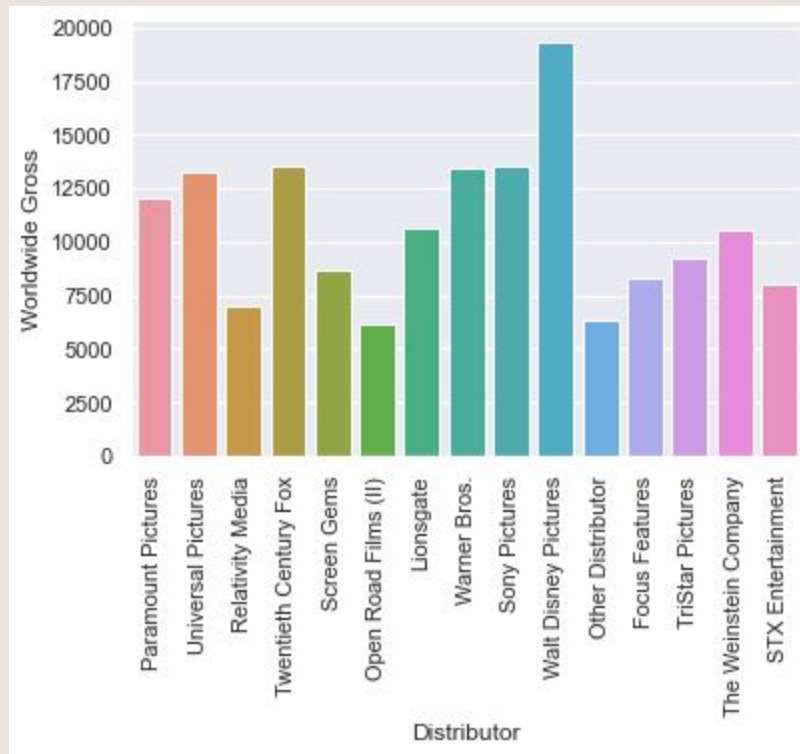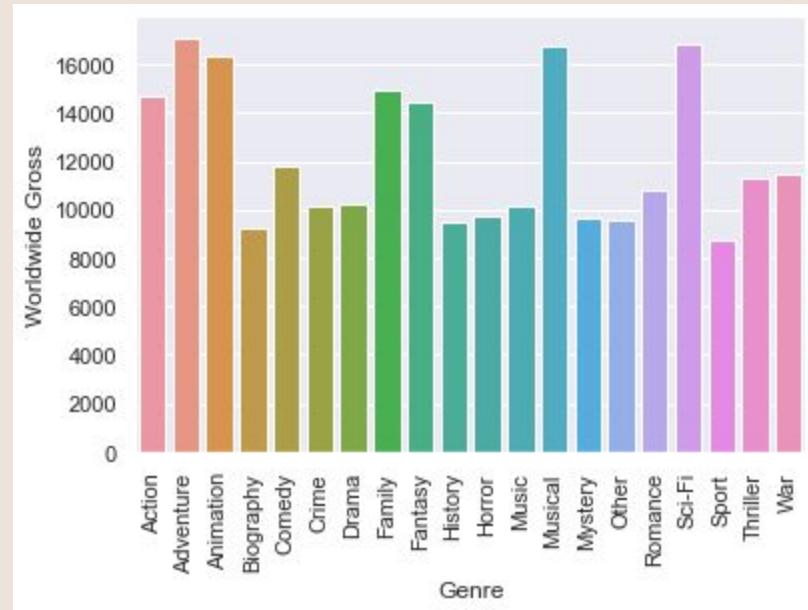- Number of Theaters
- **Distributor**
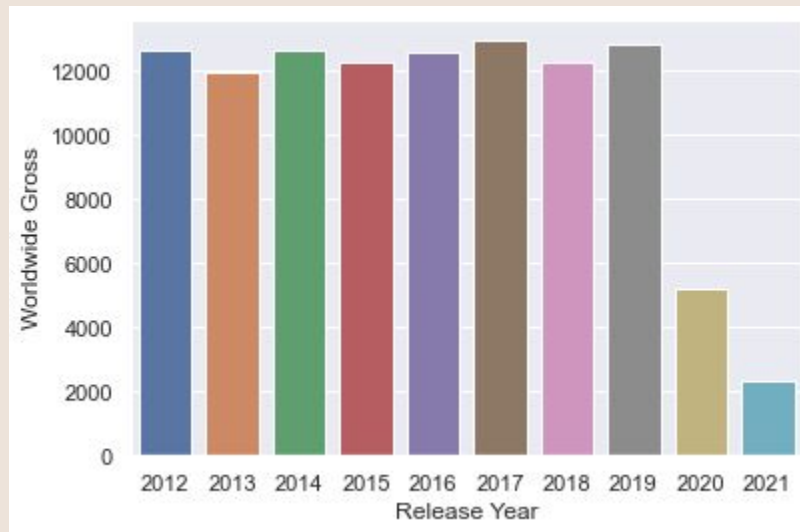- Genre
- Release Year

# Feature Investigation

- Budget
- Running Time
- Number of Theaters
- Distributor
- **Genre**
- Release Year

# Feature Investigation

- Budget
- Running Time
- Number of Theaters
- Distributor
- Genre
- **Release Year**

# Multiple Linear Regression Model

$$\sqrt{WWGross} = -7857.75 + (.000056 \cdot Budget) + (40.9309 \cdot RunTime) + (4.01770 \cdot NumTheaters) - (1596.44 \cdot OpenRdFilms) + (1617.38 \cdot Universal) + (2656.12 \cdot Disney) + (3161.74 \cdot Animation) - (1701.24 \cdot Family) - (1233.66 \cdot Fantasy) + (2178.95 \cdot Music) + (1890.15 \cdot Musical) - (2033.96 \cdot Sport) - (4644.61 \cdot Year2020) - (3862.50 \cdot Year2021)$$

# Multiple Linear Regression Model

$$\sqrt{WWGross} = -7857.75 + (.000056 \cdot Budget) + (40.9309 \cdot RunTime) + (4.01770 \cdot NumTheaters) - (1596.44 \cdot OpenRdFilms) + (1617.38 \cdot Universal) + (2656.12 \cdot Disney) + (3161.74 \cdot Animation) - (1701.24 \cdot Family) - (1233.66 \cdot Fantasy) + (2178.95 \cdot Music) + (1890.15 \cdot Musical) - (2033.96 \cdot Sport) - (4644.61 \cdot Year2020) - (3862.50 \cdot Year2021)$$

Numerical Data

# Multiple Linear Regression Model

$$\sqrt{WWGross} = -7857.75 + (.000056 \cdot Budget) + (40.9309 \cdot RunTime) + (4.01770 \cdot NumTheaters) - (1596.44 \cdot OpenRdFilms) + (1617.38 \cdot Universal) + (2656.12 \cdot Disney) + (3161.74 \cdot Animation) - (1701.24 \cdot Family) - (1233.66 \cdot Fantasy) + (2178.95 \cdot Music) + (1890.15 \cdot Musical) - (2033.96 \cdot Sport) - (4644.61 \cdot Year2020) - (3862.50 \cdot Year2021)$$

Distributors

# Multiple Linear Regression Model

$$\sqrt{WWGross} = -7857.75 + (.000056 \cdot Budget) + (40.9309 \cdot RunTime) + (4.01770 \cdot NumTheaters) - (1596.44 \cdot OpenRdFilms) + (1617.38 \cdot Universal) + (2656.12 \cdot Disney) + (3161.74 \cdot Animation) - (1701.24 \cdot Family) - (1233.66 \cdot Fantasy) + (2178.95 \cdot Music) + (1890.15 \cdot Musical) - (2033.96 \cdot Sport) - (4644.61 \cdot Year2020) - (3862.50 \cdot Year2021)$$
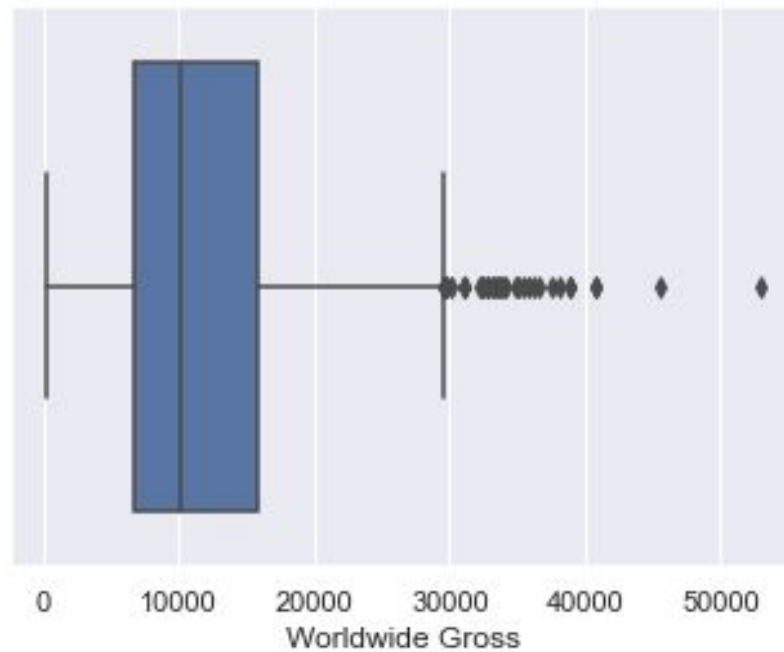
Genres

# Multiple Linear Regression Model

$$\sqrt{WWGross} = -7857.75 + (.000056 \cdot Budget) + (40.9309 \cdot RunTime) + (4.01770 \cdot NumTheaters) - (1596.44 \cdot OpenRdFilms) + (1617.38 \cdot Universal) + (2656.12 \cdot Disney) + (3161.74 \cdot Animation) - (1701.24 \cdot Family) - (1233.66 \cdot Fantasy) + (2178.95 \cdot Music) + (1890.15 \cdot Musical) - (2033.96 \cdot Sport) - (4644.61 \cdot Year2020) - (3862.50 \cdot Year2021)$$
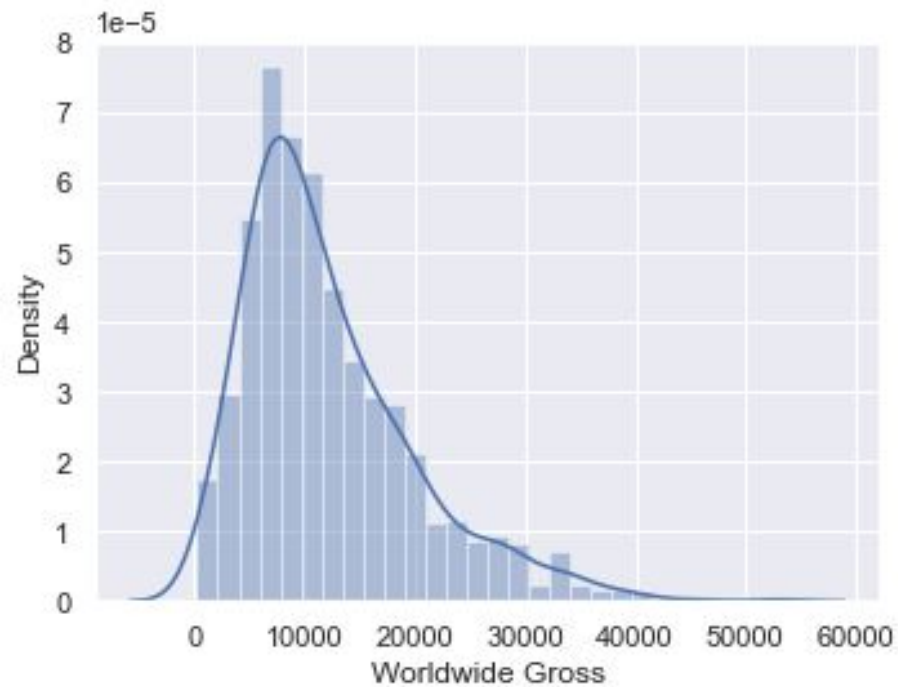
Release Years

# Conclusion

- Films with larger budget, longer running time, and showing in more theaters will likely have a larger Worldwide Gross.
- Universal Pictures and Walt Disney Pictures are more preferable distributors than Open Road Films.
- Films can benefit from including animation, music, and musicals while family, fantasy, and sport genres can be a detriment.
- Recent years have correlated with a decrease in Worldwide Gross (likely due to COVID).

# Future Directions

- Search for other important features that are likely missing from this model:
  - Production time, location, cost
  - Actors and directors
  - Pre-existing popularity or familiarity of the story
  - Extent of advertising and social media outreach
- Test out other modeling methods to find a better fit

Thank you!

# Appendix

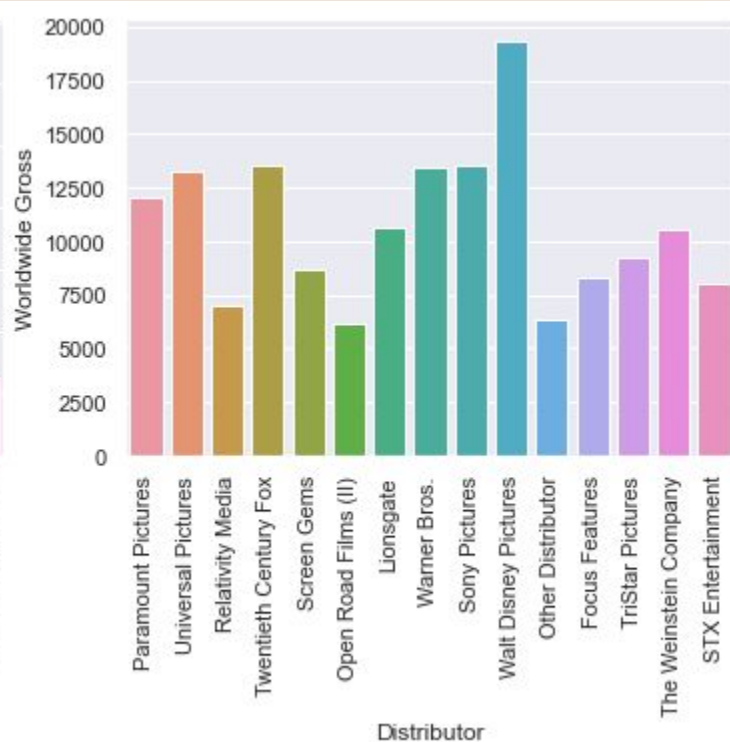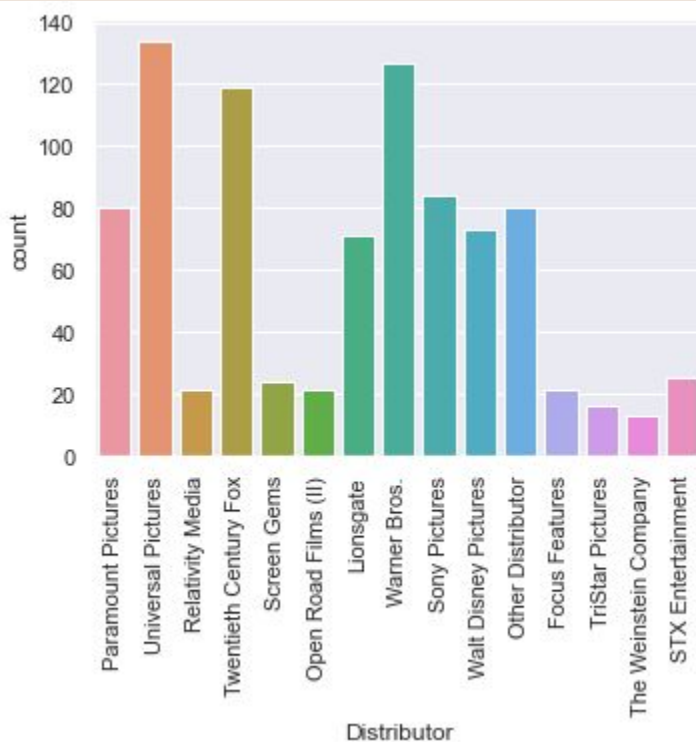# Appendix

```
Int64Index: 909 entries, 0 to 908
Data columns (total 9 columns):
 #    Column              Non-Null Count   Dtype
---   ------              --------------   -----
 0    Worldwide Gross     909 non-null     int64
 1    Distributor         909 non-null     object
 2    Budget              909 non-null     float64
 3    MPAA                909 non-null     object
 4    Running Time        909 non-null     int64
 5    Genres              909 non-null     object
 6    Number of Theaters  909 non-null     float64
 7    Release Month       909 non-null     object
 8    Release Year        909 non-null     object
dtypes: float64(2), int64(2), object(5)
memory usage: 71.0+ KB
```
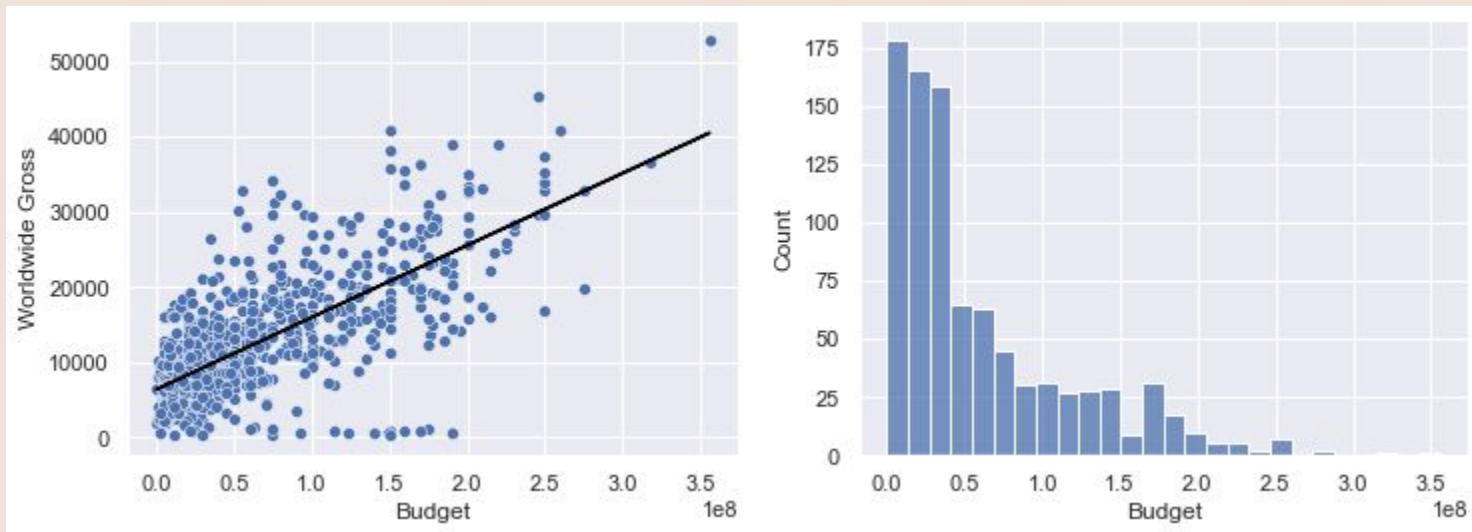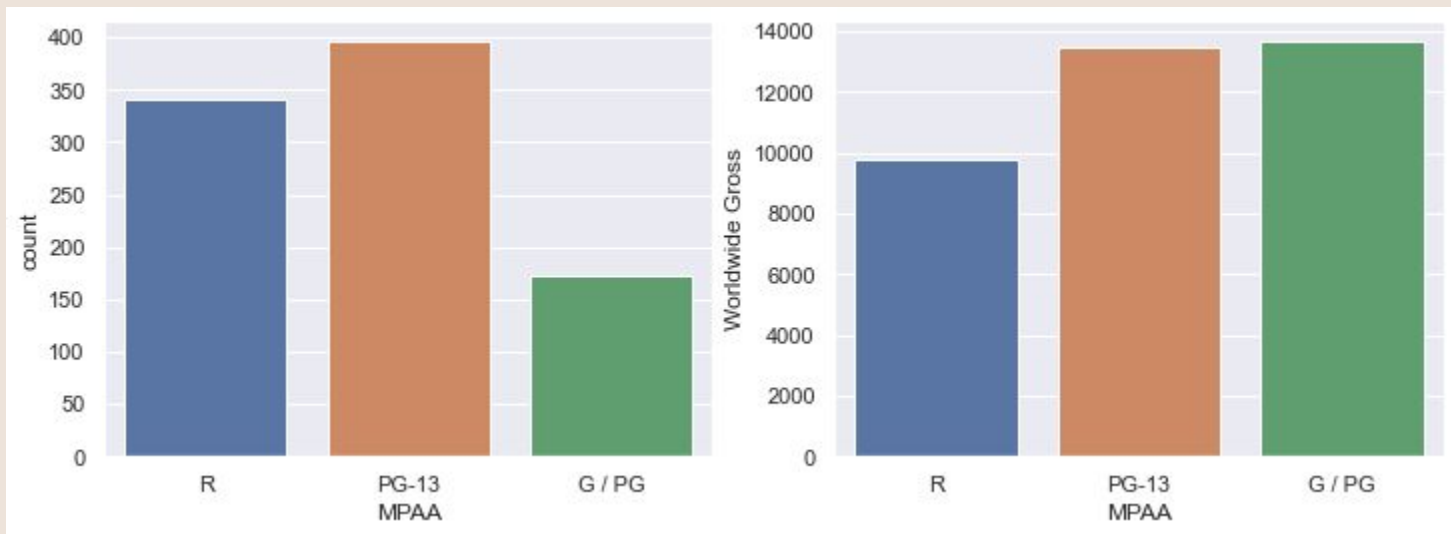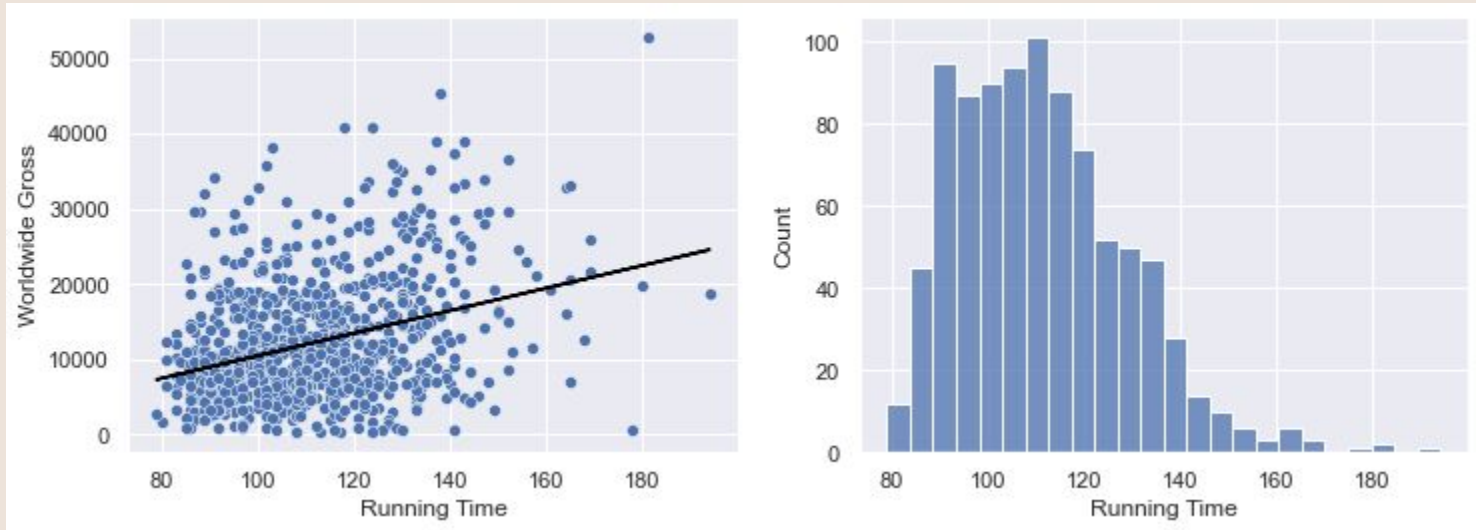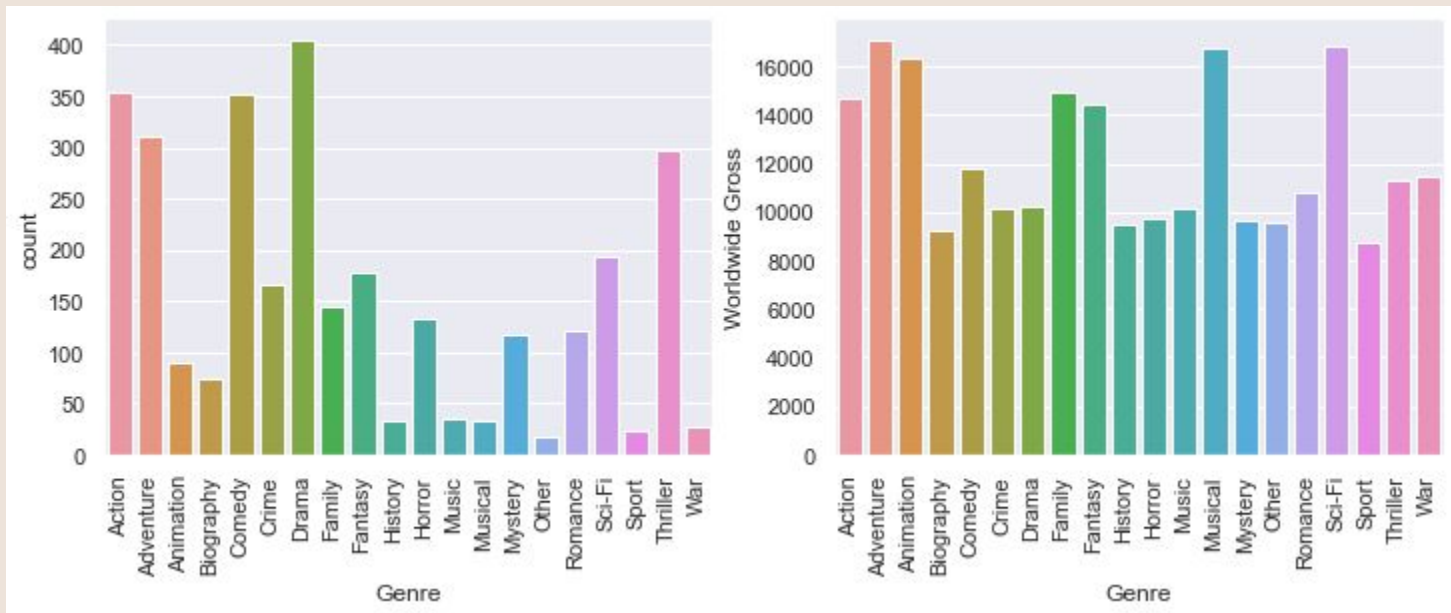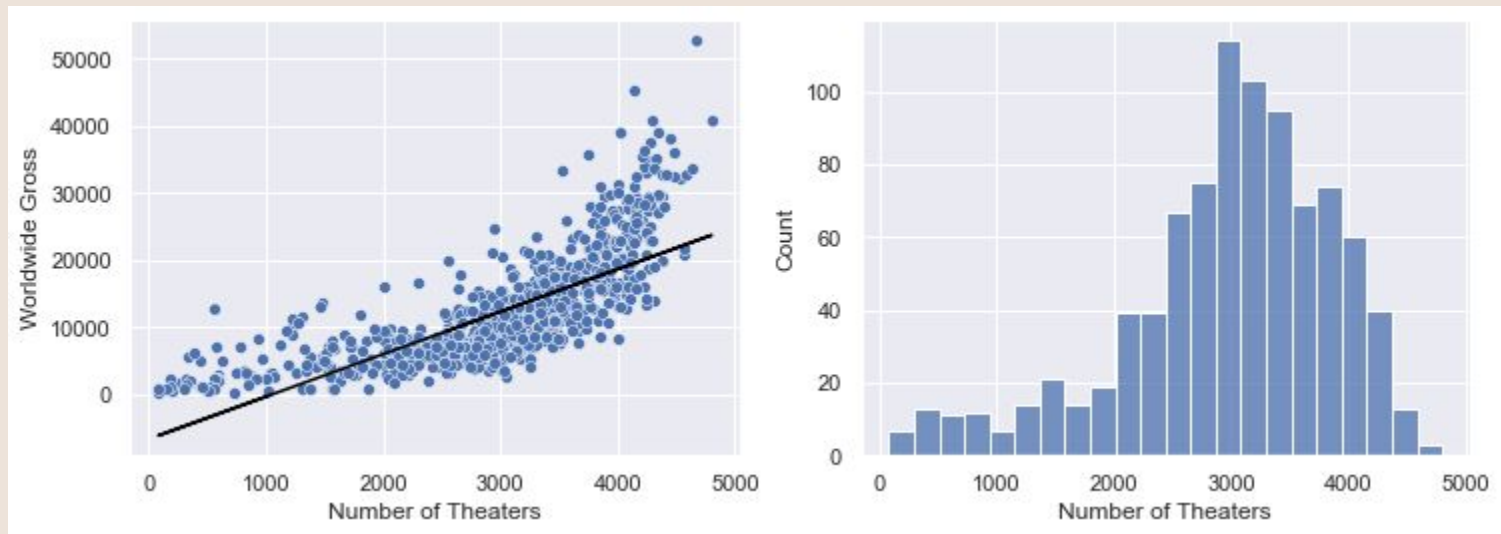
# Appendix

# Appendix

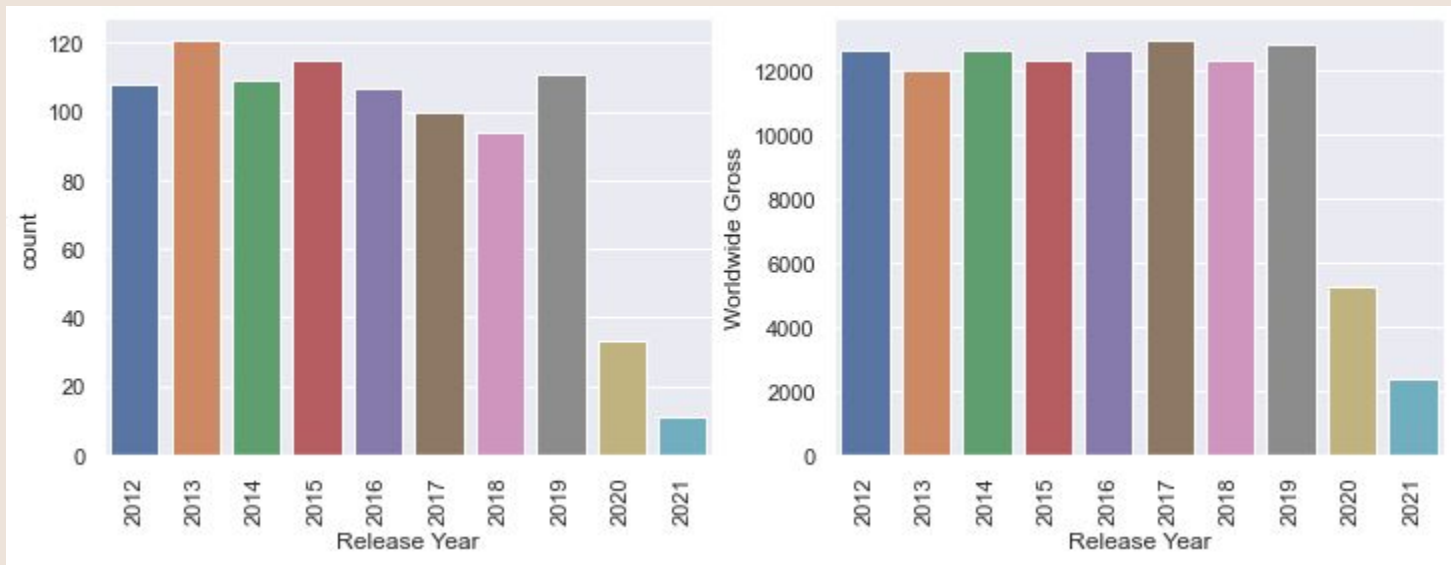# Appendix
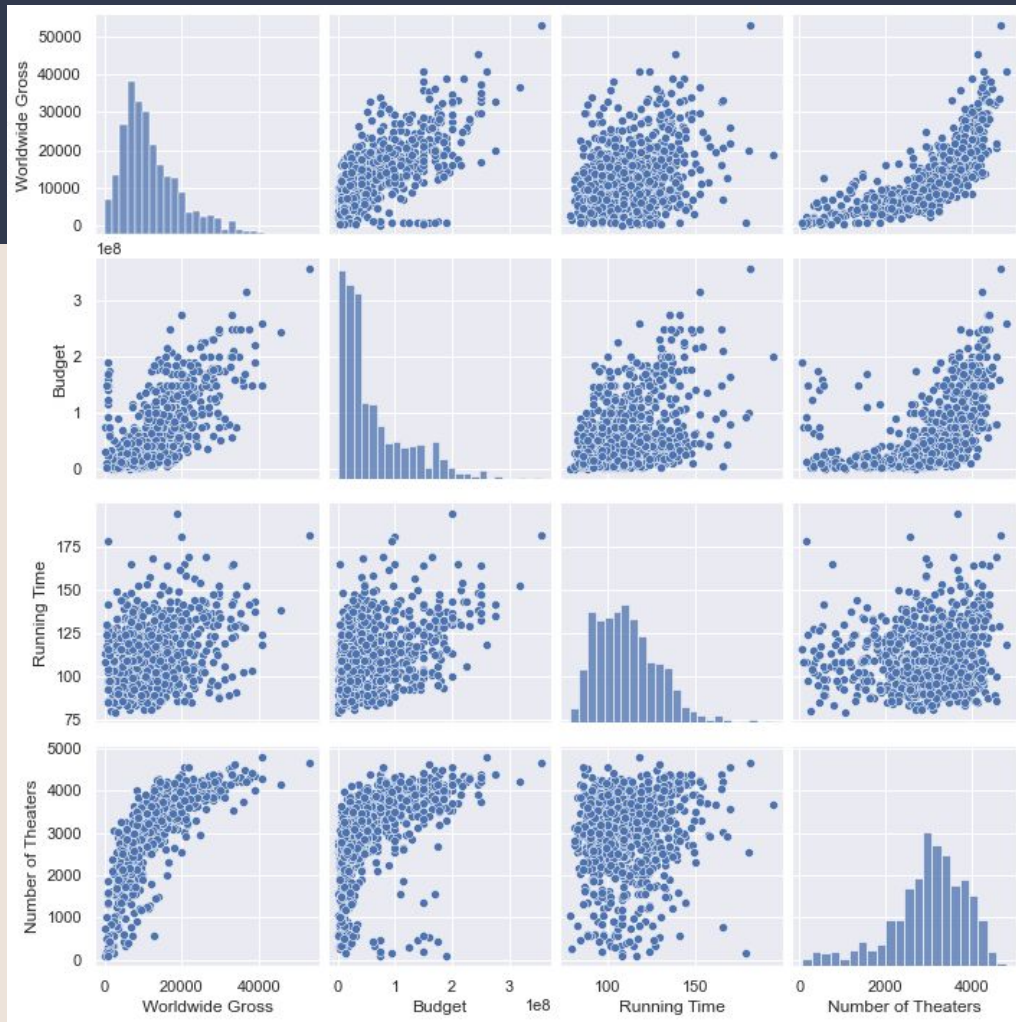
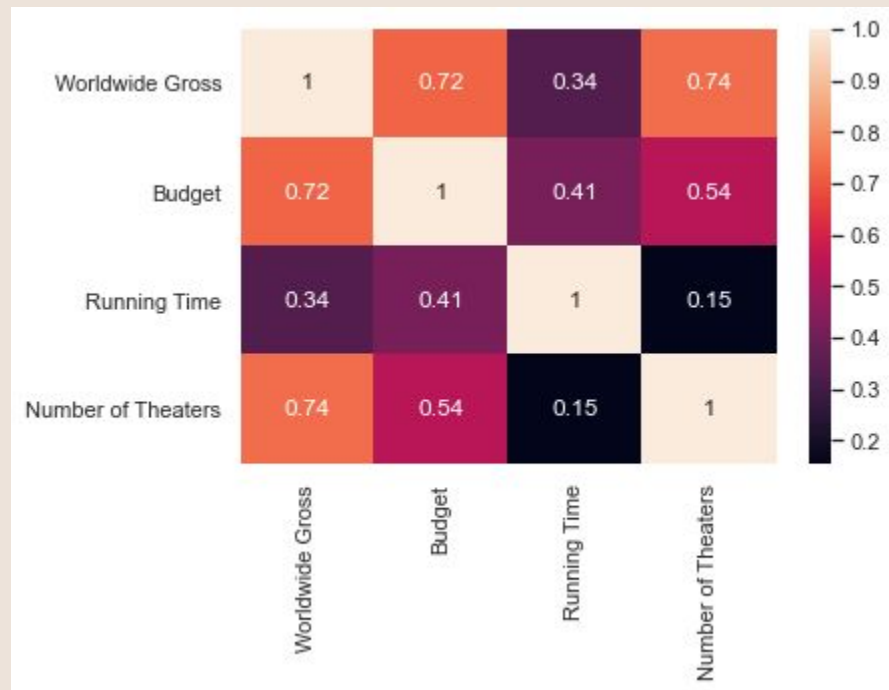# Appendix

# Appendix

# Appendix

# Appendix

# Appendix

# Appendix

# Appendix

# Appendix

```
R² training set : 0.7399702043986445
R² test set : 0.7221791414372754

Mean Squared Error : 16052250.800428031
Mean Absolute Error : 3052.5837592996995
```
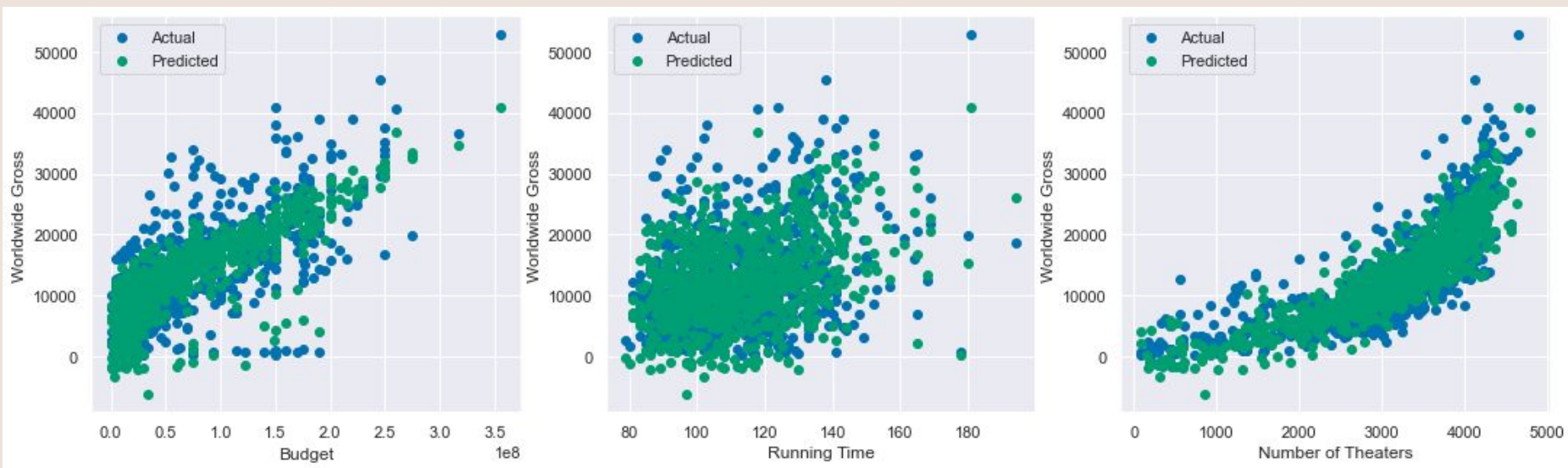
| | Features | Coefficients | Squared |
|---|---|---|---|
| 0 | Intercept | -7857.746350 | 6.174418e+07 |
| 1 | Budget | 0.000056 | 3.174432e-09 |
| 2 | Running Time | 40.930913 | 1.675340e+03 |
| 3 | Number of Theaters | 4.017704 | 1.614194e+01 |
| 4 | Open Road Films (II) | -1596.441120 | 2.548624e+06 |
| 5 | Universal Pictures | 1617.378737 | 2.615914e+06 |
| 6 | Walt Disney Pictures | 2656.119903 | 7.054973e+06 |
| 7 | Animation | 3161.738794 | 9.996592e+06 |
| 8 | Family | -1701.238735 | 2.894213e+06 |
| 9 | Fantasy | -1233.655311 | 1.521905e+06 |
| 10 | Music | 2178.947164 | 4.747811e+06 |
| 11 | Musical | 1890.150384 | 3.572668e+06 |
| 12 | Sport | -2033.955706 | 4.136976e+06 |
| 13 | 2020 | -4644.613468 | 2.157243e+07 |
| 14 | 2021 | -3862.504435 | 1.491894e+07 |

# Model Fit on Numerical Data



$R^2$ score on training set : 0.73997

$R^2$ score on test data set : 0.72218

# Appendix